## CLAIMS

1 1. A method for morphological disambiguation,
2 comprising:
3    receiving an input string;
4    morphologically analyzing the string to generate a
5 list of candidate analyses of the string, each candidate
6 analysis comprising a respective word and a linguistic
7 pattern of the word; and
8    evaluating the pattern of each of the analyses
9 against a predefined criterion in order to select one or
10 more of the analyses from the list.

1 2. A method according to claim 1, wherein receiving the
2 input string comprises receiving a word in a Semitic
3 language.

1 3. A method according to claim 2, wherein the Semitic
2 language comprises Hebrew.

1 4. A method according to claim 1, wherein the
2 linguistic pattern comprises a specification of at least
3 one characteristic of the word, selected from a set of
4 characteristics including a part of speech, prefix,
5 number, gender and person of the word.

1 5. A method according to claim 4, wherein the
2 specification of the at least one characteristic
3 comprises a specification of all of the characteristics
4 in the set.

1 6. A method according to claim 5, wherein when the base
2 word comprises a verb, the linguistic pattern further
3 comprises a designation of a tense and conjugation
4 pattern of the verb.

1 7. A method according to claim 1, wherein each of the
2 analyses has a lemma and a paradigm determined by the
3 word and the linguistic pattern thereof, and wherein
4 evaluating the pattern comprises eliminating one of the
5 analyses from the list if it has the same lemma and
6 paradigm as another of the analyses.

1 8. A method according to claim 1, wherein evaluating
2 the pattern comprises determining a relative frequency of
3 occurrence of the pattern of each of the analyses, and
4 selecting the at least one of the analyses whose
5 frequency of occurrence is above a predetermined
6 threshold.

1 9. A method according to claim 8, wherein determining
2 the relative frequency of occurrence comprises
3 morphologically analyzing a corpus of text and finding
4 the frequency of occurrence of the pattern in the corpus.

1 10. A method according to claim 9, wherein determining
2 the relative frequency of occurrence comprises storing in
3 a table the frequency of occurrence found in the corpus,
4 and looking up the pattern in the table.

1 11. A method according to claim 8, wherein selecting the
2 at least one of the analyses comprises setting the
3 threshold so as to control how many of the analyses from
4 the list are selected.

1 12. A method according to claim 8, wherein selecting the
2 at least one of the analyses comprises selecting the at
3 least one of the analyses based on the pattern thereof,
4 and substantially independently of the respective word.

1 13. A method according to claim 1, and comprising
2 searching in a corpus of text for a match to the input
3 string using the one or more selected analyses.

1 14. A method according to claim 1, and comprising
2 checking for spelling errors in the input string using
3 the one or more selected analyses.

1 15. A method for searching a corpus of text made up of
2 words, comprising:
3    morphologically analyzing the words in the corpus to
4 generate, for each of at least some of the words, a list
5 of candidate analyses, each candidate analysis comprising
6 a respective lemma and a linguistic pattern relating the
7 lemma to the analyzed word;
8    evaluating the pattern of each of the analyses
9 against a predefined criterion in order to select one or
10 more of the analyses from the list for each of the
11 analyzed words;
12    entering the lemmas of the selected analyses in an
13 index of the corpus; and
14    applying a search query to the index.

1 16. A method according to claim 15, wherein applying the
2 search query comprises:
3    receiving an input text string;
4    morphologically analyzing and disambiguating the
5 string to generate one or more search lemmas for the
6 string; and
7    comparing the search lemmas to the index.

1 17. A method according to claim 15, wherein the words in
2 the corpus comprise words in a Semitic language.

1 18. A method according to claim 17, wherein the Semitic
2 language comprises Hebrew.

19. A method according to claim 15, wherein the linguistic pattern comprises a specification of at least one characteristic of the word, selected from a set of characteristics including a part of speech, prefix, number, gender and person of the word.

20. A method according to claim 15, wherein evaluating the pattern comprises determining a relative frequency of occurrence of the pattern of each of the analyses, and selecting the at least one of the analyses whose frequency of occurrence is above a predetermined threshold.

21. A method according to claim 20, wherein selecting the at least one of the analyses comprises selecting the at least one of the analyses based on the pattern thereof, and substantially independently of the respective word.

22. A computer software product, comprising a computer-readable medium in which program instructions are stored, which instructions, when read by a computer, cause the computer to morphologically analyze an input string to generate a list of candidate analyses of the string, each candidate analysis comprising a respective word and a linguistic pattern of the word, and to evaluate the pattern of each of the analyses against a predefined criterion in order to select one or more of the analyses from the list.

23. A product according to claim 22, wherein the input string comprises a word in a Semitic language.

24. A product according to claim 23, wherein the Semitic language comprises Hebrew.

1 25. A product according to claim 22, wherein the
2 linguistic pattern comprises a specification of at least
3 one characteristic of the word, selected from a set of
4 characteristics including a part of speech, prefix,
5 number, gender and person of the word.

1 26. A product according to claim 22, wherein the
2 instructions further cause the computer to search in a
3 corpus of text for a match to the input string using the
4 one or more selected analyses.

1 27. A computer software product, comprising a
2 computer-readable medium in which program instructions
3 are stored, which instructions, when read by a computer,
4 cause the computer to morphologically analyze the words
5 in the corpus to generate, for each of at least some of
6 the words, a list of candidate analyses, each candidate
7 analysis comprising a respective lemma and a linguistic
8 pattern relating the lemma to the analyzed word, to
9 evaluate the pattern of each of the analyses against a
10 predefined criterion in order to select one or more of
11 the analyses from the list for each of the analyzed
12 words, to enter the lemmas of the selected analyses in an
13 index of the corpus, and to apply a search query to the
14 index.

1 28. A product according to claim 27, wherein the
2 instructions further cause the computer to receive an
3 input text string, to morphologically analyze and
4 disambiguate the string to generate one or more search
5 lemmas for the string, and to compare the search lemmas
6 to the index.

1 29. Apparatus for morphological disambiguation,
2 comprising a linguistic processor, which is adapted to

3   receive an input string, to morphologically analyze the
4   string to generate a list of candidate analyses of the
5   string, each candidate analysis comprising a respective
6   word and a linguistic pattern of the word, and to
7   evaluate the pattern of each of the analyses against a
8   predefined criterion in order to select one or more of
9   the analyses from the list.

1   30.  Apparatus according to claim 29, wherein the input
2   string comprises a word in a Semitic language.

1   31.  Apparatus according to claim 30, wherein the Semitic
2   language comprises Hebrew.

1   32.  Apparatus according to claim 29, wherein the
2   linguistic pattern comprises a specification of at least
3   one characteristic of the word, selected from a set of
4   characteristics including a part of speech, prefix,
5   number, gender and person of the word.

1   33.  Apparatus according to claim 29, wherein the
2   processor is further adapted to search in a corpus of
3   text for a match to the input string using the one or
4   more selected analyses.

1   34.  Apparatus for searching a corpus of text made up of
2   words, comprising a linguistic processor, which is
3   adapted to morphologically analyze the words in the
4   corpus to generate, for each of at least some of the
5   words, a list of candidate analyses, each candidate
6   analysis comprising a respective lemma and a linguistic
7   pattern relating the lemma to the analyzed word, to
8   evaluate the pattern of each of the analyses against a
9   predefined criterion in order to select one or more of
10  the analyses from the list for each of the analyzed
11  words, to enter the lemmas of the selected analyses in an

12  index of the corpus, and to apply a search query to the
13  index.

1  35.  Apparatus  according  to  claim  34,  wherein  the
2  processor  is  further  adapted  to  receive  an  input  text
3  string,  to  morphologically  analyze  and  disambiguate  the
4  string  to  generate  one  or  more  search  lemmas  for  the
5  string,  and  to  compare  the  search  lemmas  to  the  index.